

On Indeterminate Strings Matching

Paweł Gawrychowski

Institute of Computer Science, University of Wrocław, Poland

Samah Ghazawi¹

Department of Computer Science, University of Haifa, Haifa, Israel

idrees.samah@gmail.com

Gad M. Landau

Department of Computer Science, University of Haifa, Haifa, Israel

Department of Computer Science and Engineering, NYU Tandon School of Engineering, Brooklyn, USA

Abstract

Given two indeterminate equal-length strings p and t with a set of characters per position in both strings, we obtain a determinate string p_w from p and a determinate string t_w from t by choosing one character per position. Then, we say that p and t match when p_w and t_w match for some choice of the characters. While the most standard notion of a match for determinate strings is that they are simply identical, in certain applications it is more appropriate to use other definitions, with the prime examples being parameterized matching, order-preserving matching, and the recently introduced Cartesian tree matching. We provide a systematic study of the complexity of string matching for indeterminate equal-length strings, for different notions of matching. We use n to denote the length of both strings, and r to be an upperbound on the number of uncertain characters per position. First, we provide the first polynomial time algorithm for the Cartesian tree version that runs in deterministic $\mathcal{O}(n \log^2 n)$ and expected $\mathcal{O}(n \log n \log \log n)$ time using $\mathcal{O}(n \log n)$ space, for constant r . Second, we establish NP-hardness of the order-preserving version for $r = 2$, thus solving a question explicitly stated by Henriques et al. [CPM 2018], who showed hardness for $r = 3$. Third, we establish NP-hardness of the parameterized version for $r = 2$. As both parameterized and order-preserving indeterminate matching reduce to the standard determinate matching for $r = 1$, this provides a complete classification for these three variants.

2012 ACM Subject Classification Theory of computation → Pattern matching

Keywords and phrases pattern matching, indeterminate strings, Cartesian trees, order-preserving matching, parameterized matching

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

Funding *Samah Ghazawi*: Partially supported by the Israel Science Foundation (ISF) grant 1475/18.

Gad M. Landau: Partially supported by the Israel Science Foundation (ISF) grant 1475/18, and the United States-Israel Binational Science Foundation (BSF) grant No. 2018141.

¹ Corresponding author.



© Paweł Gawrychowski, Samah Ghazawi and Gad M. Landau;
licensed under Creative Commons License CC-BY

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:12

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

1 Introduction

String matching, in the sense of comparing two equal-length strings, is one of the fundamental problems in computer science with multiple practical applications. While exact matching is trivial to solve in optimal linear time by comparing the strings character-by-character, for many of the applications it seems more appropriate to work with some kind of approximate matching. Prime examples include string matching with swaps [2], parameterized string matching [6], string matching with gaps [9], jumbled string matching [10], string matching with don't cares [29], and edit distance [32]. In all of such problems, one needs to first precisely define when do two strings match.

Parameterized matching is a classical notion motivated by finding identical sections of code [3, 4, 5, 6, 19, 34]. Formally, two strings p and t of length n are a parameterized match when for every $i, j \in \{1, \dots, n\}$, $p[i] = p[j]$ iff $t[i] = t[j]$. This is denoted by $p \sim_{=} t$.

Order-preserving matching is a more recent but already well-studied notion motivated by stock price analysis and musical melody matching [11, 16, 17, 25, 26]. Formally, two strings p and t of length n are a order-preserving match when for every $i, j \in \{1, \dots, n\}$, $p[i] \leq p[j]$ iff $t[i] \leq t[j]$. This is denoted by $p \sim_{\leq} t$.

Very recently, a different notion called Cartesian tree matching has been proposed [28]. The Cartesian tree of a given string p ($CT(p)$), first defined in [31], is constructed according to the following rules:

- If p is an empty string, $CT(p)$ is an empty tree.
- If $p[1..n]$ is not empty and $p[i]$ is the leftmost minimum value in p , $CT(p)$ is the tree with $p[i]$ being the root, $CT(p[1..i-1])$ the left subtree, and $CT(p[i+1..n])$ the right subtree.

Even though the most well-known applications of Cartesian trees are probably in designing space-efficient structures for finding the minimum in a range, they can be also used to compare strings. Similarly to order-preserving matching, this notion is motivated by applications concerned with time-series data such as stock price analysis, and has gained considerable attention during the last year [7, 18, 30]. Formally, two strings p and t of equal length n are a Cartesian tree match when their Cartesian trees $CT(p)$ and $CT(t)$ are identical. This is denoted by $p \sim_C t$.

We consider the complexity of string matching for indeterminate strings defined as follows.

► **Definition 1.** An indeterminate string is a sequence of sets of characters $p[1]p[2]...p[n]$, where $p[i] \subseteq \mathbb{N}$. Each position is specified by writing $p[i] = x_1|...|x_r$, such that $x_\ell \in \mathbb{N}$, which means that we can choose $p[i]$ to be any x_ℓ .

Indeterminate strings were studied earlier, among others, covering problems for indeterminate strings [1, 14] and indeterminate strings in graph theory [20, 12, 27]. Indeterminate string matching was investigated lately from different angles [8, 13, 15, 22, 23, 24]. It provides a convenient formalism for compactly capturing situations in which there are some uncertainties concerning characters at some positions. Indeed, an indeterminate string p of length n describes r^n determinate strings. We write \tilde{p} to denote the set of all such strings, and p_w when referring to a single determinate string described by p .

First, we consider the complexity of Cartesian tree matching for indeterminate strings defined as follows.

Problem: CARTESIAN TREE MATCHING OF INDETERMINATE STRINGS (CTMIS)

Input: Two indeterminate strings p and t of length n with up to r of uncertain characters per position.

Output: Are there determinate strings $p_w \in \tilde{p}$ and $t_w \in \tilde{t}$ such that p_w Cartesian tree matches t_w ?

A naive solution to the CTMIS would be to apply the solution of [28] to each $t_w \in \tilde{t}$ and $p_w \in \tilde{p}$ in $\mathcal{O}(n^2 r^n)$ time. In Section 2 we provide the first polynomial algorithm for this problem that works in $\mathcal{O}(n \log^2 n)$ time and $\mathcal{O}(n \log n)$ space, assuming that r is constant. Additionally, in the Word RAM model of computation we further improve the time complexity to expected $\mathcal{O}(n \log n \log \log n)$.

► **Example 2.** Consider the following indeterminate strings:

$$p = (2|4|7, 2|5|6, 1|4|8, 4|7|8, 3|10|16)$$

$$t = (2|7|10, 5|20|31, 10|17|25, 0|9|11, 1|8|18).$$

$p_w = (7, 2, 8, 4, 16)$ and $t_w = (10, 5, 17, 9, 18)$ define the same Cartesian tree, see Figure 1. Therefore, we say that $p \sim_C t$. Note that p and t define other matching or non-matching Cartesian trees.



■ **Figure 1** The Cartesian trees of $p_w = (7, 2, 8, 4, 16)$ and $t_w = (10, 5, 17, 9, 18)$.

Second, we consider the complexity of order-preserving matching for indeterminate strings defined as follows.

Problem: ORDER-PRESERVING MATCHING OF INDETERMINATE STRINGS (OPMIS)

Input: Two indeterminate strings p and t of length n with up to r uncertain characters per position.

Output: Are there determinate strings $p_w \in \tilde{p}$ and $t_w \in \tilde{t}$ such that p_w order-preserving matches t_w ?

Henriques et al. [21] proved that OPMIS is NP-hard for $r = 3$. As for $r = 1$ there is a simple linear-time algorithm, this left $r = 2$ as the only open case (CPM version of the paper [21] claims a polynomial time algorithm for this case, but this has been clarified in the arXiv version [13]). In Section 4 we provide a different reduction that establishes NP-hardness of OPMIS already for $r = 2$, thus fully resolving the complexity of this problem and answering an open question explicitly stated by Costa et al. In contrast with the previous work, our reduction exploits the order between elements instead of just their equality, and is more involved.

Third, we consider the complexity of parameterized matching for indeterminate strings defined as follows.

Problem: PARAMETERIZED MATCHING OF INDETERMINATE STRINGS (PMIS)

Input: Two indeterminate strings p and t of length n with up to r uncertain characters per position.

Output: Are there determinate strings $p_w \in \tilde{p}$ and $t_w \in \tilde{t}$ such that p_w parameterized matches t_w ?

NP-hardness proof by Henriques et al. [21] implicitly shows hardness of PMIS for $r = 3$. This, again, leaves $r = 2$ as the only open case. In Section 5 we provide a reduction that establishes NP-hardness of PMIS for $r = 2$.

2 CTMIS in $\mathcal{O}(n^3)$ Time and $\mathcal{O}(n^2)$ Space

In this section, we describe a warm-up solution for the CTMIS problem. The input is two equal-length indeterminate strings p and t with two uncertain characters per position, and the output is whether $p \sim_C t$ or not. The solution can be generalized to any constant value of r in a straightforward manner. We will assume that both p and t consists of distinct values, which can be always ensured by an appropriate perturbation.

First, note that for each index i , we have $p[i] = x_i|x'_i$ and $t[i] = y_i|y'_i$, hence each i defines 4 pairs $\{(x_i, y_i), (x_i, y'_i), (x'_i, y_i), (x'_i, y'_i)\}$ called thresholds. The main idea of the algorithm is to determine for each index i and a threshold (x_i, y_i) :

1. for which indices k we have $p[k, i] \sim_C t[k, i]$ with the roots x_i and y_i , respectively.
2. for which indices j we have $p[i, j] \sim_C t[i, j]$ with the roots x_i and y_i , respectively.

Consider an interval $[k, i]$, the reasoning for an interval $[i, j]$ is similar. We have $p[k, i] \sim_C t[k, i]$ with the roots x_i and y_i iff there exists an index ℓ and a threshold (x_ℓ, y_ℓ) where $k \leq \ell \leq i - 1$, $x_i < x_\ell$ and $y_i < y_\ell$ such that $p[k, \ell] \sim_C t[k, \ell]$ and $p[\ell, i - 1] \sim_C t[\ell, i - 1]$.

We process all possible intervals $[k, i]$ and $[i, j]$ in an increasing order of their lengths using dynamic programming. For each index i and threshold (x_i, y_i) we compute the answer for all left intervals $[k, i - 1]$ and all right intervals $[i + 1, j]$, see Figure 2. We define two types of states and associate a boolean value with each of them as follows:

Left states $L_{k,i}(x_i, y_i) = \text{true}$ iff $p[k, i] \sim_C t[k, i]$ with the roots x_i and y_i , respectively.

Right states $R_{i,j}(x_i, y_i) = \text{true}$ iff $p[i, j] \sim_C t[i, j]$ with the roots x_i and y_i , respectively.

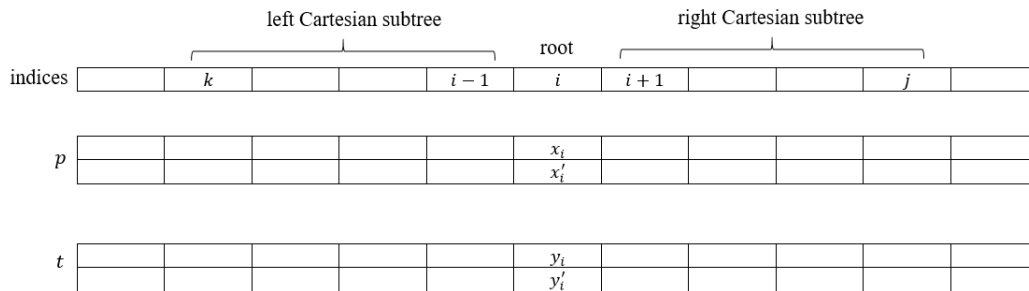


Figure 2 An interval $[k, j]$ with the root at index i , defining left and right Cartesian subtrees.

► **Example 3.** Let $p = (4|7, 2|6, 1|8, 3|20, 10|16)$ and $t = (2|10, 20|31, 10|17, 0|11, 8|18)$. Considering index 3 as a possible root for the Cartesian tree in both strings, the thresholds defined by index 3 are $(x_3, y_3) \in \{(1, 10), (8, 10), (1, 17), (8, 17)\}$. The right states are $R_{3,4}(x_3, y_3)$ and

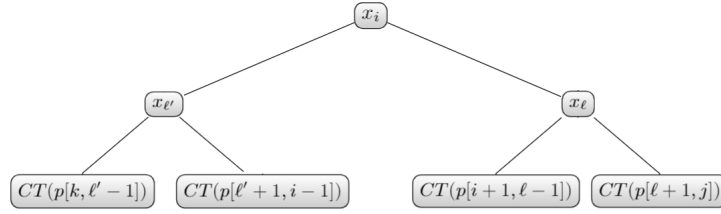
23:4 On Indeterminate Strings Matching

129 $R_{3,5}(x_3, y_3)$. The left states are $L_{2,3}(x_3, y_3)$ and $L_{1,3}(x_3, y_3)$. Some of their corresponding
130 boolean values are as follows:

- 131 1. $R_{3,4}(1, 10) = \text{true}$ for $p[3, 4] = (1, 3)$ and $t[3, 4] = (10, 11)$.
- 132 2. $R_{3,4}(8, 10) = \text{true}$ for $p[3, 4] = (8, 20)$ and $t[3, 4] = (10, 11)$.
- 133 3. $L_{2,3}(1, 10) = \text{true}$ for $p[2, 3] = (6, 1)$ and $t[2, 3] = (31, 10)$.
- 134 4. $L_{1,3}(1, 10) = \text{true}$ for $p[1, 3] = (4, 6, 1)$ and $t[2, 3] = (2, 31, 10)$.

135 From the definition of a Cartesian tree we directly obtain the following proposition
136 illustrated in Figure 3.

137 ► **Proposition 4.** (a) $R_{i,j}(x_i, y_i) = \text{true}$ iff $\exists \ell \in [i + 1, j]$ such that $R_{\ell,j}(x_\ell, y_\ell) = \text{true}$ and
138 $L_{i+1,\ell}(x_\ell, y_\ell) = \text{true}$ where $x_\ell > x_i$ and $y_\ell > y_i$. (b) $L_{k,i}(x_i, y_i) = \text{true}$ iff $\exists \ell' \in [k, i - 1]$
139 such that $R_{\ell',i-1}(x_{\ell'}, y_{\ell'}) = \text{true}$ and $L_{k,\ell'}(x_{\ell'}, y_{\ell'}) = \text{true}$ where $x_{\ell'} > x_i$ and $y_{\ell'} > y_i$.



■ **Figure 3** The Cartesian tree of $p[k, j]$ with i position as the root. The left Cartesian subtree corresponds for the left state $L_{k,i}(x_i, y_i)$. The right Cartesian subtree corresponds for the right state $R_{i,j}(x_i, y_i)$. Note that, $t[k, j]$ matches the tree above with values y_i, y_ℓ and $y_{\ell'}$.

140 Recall that we apply dynamic programming in an increasing order of the lengths of
141 the intervals. Therefore, the states $R_{\ell,j}(x_\ell, y_\ell)$ and $L_{i+1,\ell}(x_\ell, y_\ell)$ from Proposition 4(a) are
142 computed before the state $R_{i,j}(x_i, y_i)$. Similarly, the states $R_{\ell',i-1}(x_{\ell'}, y_{\ell'})$ and $L_{k,\ell'}(x_{\ell'}, y_{\ell'})$
143 are computed before the state $L_{k,i}(x_i, y_i)$. Therefore, for every interval we can simply consider
144 all relevant ℓ and ℓ' , access their corresponding states, and update the answer. Finally, after
145 having processed all the intervals, we conclude that $p \sim_C t$ iff there exists an index i and a
146 threshold (x_i, y_i) such that $L_{1,i}(x_i, y_i) = \text{true}$ and $R_{i,n}(x_i, y_i) = \text{true}$.

147 ► **Example 5.** Let $p = (4|7, 2|6, 1|8, 3|20, 10|16)$ and $t = (2|10, 20|31, 10|17, 0|11, 8|18)$ as
148 in the previous example above. The states: $L_{1,3}(1, 10) = \text{true}$ for $p[1, 3] = (4, 6, 1)$ and
149 $t[2, 3] = (2, 31, 10)$, and $R_{3,5}(1, 10) = \text{true}$ for $p[3, 5] = (1, 3, 16)$ and $t[3, 5] = (10, 11, 18)$.
150 Hence, $p \sim_C t$ with the roots 1 and 10 respectively.

151 **Time complexity.** For each state $L_{k,i}(x_i, y_i)$ and $R_{i,j}(x_i, y_i)$ we consider $\mathcal{O}(n)$ relevant
152 indices ℓ and ℓ' , respectively. Each such index is processed in constant time, thus the overall
153 time complexity is $\mathcal{O}(n^3)$. The space complexity is bounded by the number of states processed
154 in the dynamic programming, which is $\mathcal{O}(n^2)$.

155 3 CTMIS $\mathcal{O}(n \log^2 n)$ Time and $\mathcal{O}(n \log n)$ Space

156 In this section we present an efficient solution for the CTMIS problem that builds on the
157 slower algorithm presented in the previous section.

158 The input is two equal-length indeterminate strings p and t with 2 uncertain characters
159 per position, and the output is whether $p \sim_C t$, or not. The solution can be generalized to

any constant value of r in a straightforward manner. The main idea of the algorithm is to find, for index i and a threshold (x_i, y_i) , the largest matching Cartesian trees with the root in both trees being x_i and y_i at position i , respectively. As in the previous algorithm, we consider each index i and a threshold (x_i, y_i) separately. However, now instead of computing the answer for all intervals $[k, i]$ and $[i, j]$ we use the following definition.

► **Definition 6.** For an index i and a threshold (x_i, y_i) :

- $\min_L(x_i, y_i)$ denotes the smallest index such that $p[\min_L(x_i, y_i), i] \sim_C t[\min_L(x_i, y_i), i]$,
- $\max_R(x_i, y_i)$ denotes the largest index such that $p[i, \max_R(x_i, y_i)] \sim_C t[i, \max_R(x_i, y_i)]$,

with the root in both trees being x_i and y_i at position i , respectively.

Computing $\min_L(x_i, y_i)$ and $\max_R(x_i, y_i)$ fully describes the situation, as the above definition together with the definition of a Cartesian tree matching directly imply the following:

- $p[\ell, i] \sim_C t[\ell, i]$ iff $\min_L(x_i, y_i) \leq \ell \leq i$.
- $p[i, r] \sim_C t[i, r]$ iff $i \leq r \leq \max_R(x_i, y_i)$.

We also note that $p[\min_L(x_i, y_i), \max_R(x_i, y_i)] \sim_C t[\min_L(x_i, y_i), \max_R(x_i, y_i)]$ due to a Cartesian tree with the root in both trees being x_i and y_i at position i , respectively. Consequently, $p \sim_C t$ iff there exists an index i and a threshold (x_i, y_i) such that $\min_L(x_i, y_i) = 1$ and $\max_R(x_i, y_i) = n$. Thus, in the remaining part of this section we focus on efficiently computing the values of \min_L and \max_R .

Our algorithm is based on the following observation. Consider an index i and a threshold (x_i, y_i) , and assume that $\min_L(x_i, y_i)$ and $\max_R(x_i, y_i)$ have been already computed. Then, the following holds:

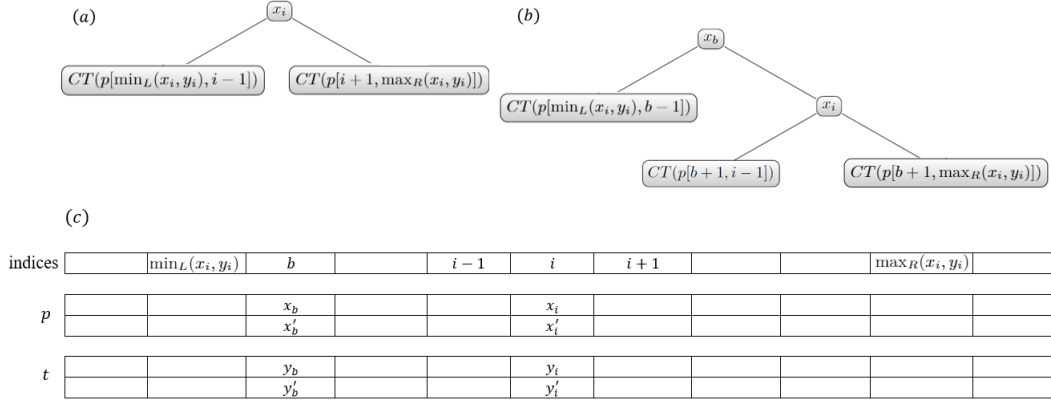
1. for any index $b \in [\min_L(x_i, y_i) - 1, \max_R(x_i, y_i)]$ and a threshold (x_b, y_b) such that $x_b < x_i$ and $y_b < y_i$, the index $\max_R(x_i, y_i)$ is a potential candidate for $\max_R(x_b, y_b)$.
2. for any index $b \in [\min_L(x_i, y_i), \max_R(x_i, y_i) + 1]$ and a threshold (x_b, y_b) such that $x_b < x_i$ and $y_b < y_i$, the index $\min_L(x_i, y_i)$ is a potential candidate for $\min_L(x_b, y_b)$.

Each index b and a threshold (x_b, y_b) might be considered for several indices i and thresholds (x_i, y_i) in the above statement, hence we might have several potential candidates for $\min_L(x_b, y_b)$ and $\max_R(x_b, y_b)$. By the definition of a Cartesian tree, one of these potential candidates corresponds to the sought $\min_L(x_b, y_b)$ and $\max_R(x_b, y_b)$ as defined above.

The high-level description of the algorithm is as follows. We iterate over all indices i and thresholds (x_i, y_i) in a specific order that will be precisely defined later. For each index i and a threshold (x_i, y_i) we aim to:

- Step 1** Compute efficiently the indices $\min_L(x_i, y_i)$ and $\max_R(x_i, y_i)$, where $\min_L(x_i, y_i)$ equals the minimum among the potential candidates for $\min_L(x_i, y_i)$ and $\max_R(x_i, y_i)$ equals the maximum among the potential candidates for $\max_R(x_i, y_i)$.
- Step 2** Update for all thresholds (x_b, y_b) at indices b such that $x_b < x_i$ and $y_b < y_i$ where $\min_L(x_i, y_i) \leq b \leq \max_R(x_i, y_i) + 1$ as a potential candidate $\min_L(x_i, y_i)$. Additionally, where $\min_L(x_i, y_i) - 1 \leq b \leq \max_R(x_i, y_i)$ as a potential candidate $\max_R(x_i, y_i)$.

We need to ensure that, for any index i and a threshold (x_i, y_i) , and any index b and a threshold (x_b, y_b) such that $x_b < x_i$ and $y_b < y_i$, $\min_L(x_i, y_i)$ and $\max_R(x_i, y_i)$ are already



■ **Figure 4** (a) The Cartesian tree of $p[\min_L(x_i, y_i), \max_R(x_i, y_i)]$ with x_i as a root. (b) The Cartesian tree of $p[\min_L(x_i, y_i), \max_R(x_i, y_i)]$ with x_b as the root after changing x'_b to x_b . The Cartesian trees of $t[\min_L(x_i, y_i), \max_R(x_i, y_i)]$ matches the trees above with the values y_b and y_i . (c) The given strings p and t with the proper intervals defining the Cartesian trees (a) and (b) above.

202 computed when we are considering threshold (x_b, y_b) at index b . This will be guaranteed by
 203 the algorithm as explained below.

204 The algorithm considers all indices i and thresholds (x_i, y_i) in the reverse lexicographical
 205 order, that is, the decreasing order of x_i and, if there is a tie, the decreasing order of y_i . Before
 206 we explain how to implement **Step 1** and **Step 2** efficiently, we define the necessary data
 207 structures. We maintain a balanced binary search tree T_y on the values of y_i , and identify y_i
 208 with its corresponding node of T_y . In each node u of T_y we have its associated secondary
 209 trees $T_{\min}(u)$ and $T_{\max}(u)$. Each $T_{\min}(u)$ and $T_{\max}(u)$ stores a collection of intervals $[\ell, r]$.
 210 The update adds a new interval $[\ell, r]$ to the collection. The query in $T_{\min}(u)$ for i finds
 211 the smallest ℓ such that $[\ell, r]$ containing i belongs to the collection, while the query in
 212 $T_{\max}(u)$ finds the largest r . By symmetry, it is enough to explain how to implement $T_{\min}(u)$.
 213 We maintain the following invariant: there are no two intervals $[\ell, r]$ and $[\ell', r']$ such that
 214 $[\ell, r] \subseteq [\ell', r']$. Clearly, such $[\ell, r]$ is not an answer to any query. Note that this implies that
 215 if we sort all the remaining intervals $[\ell_1, r_1], [\ell_2, r_2], \dots, [\ell_s, r_s]$ so that $\ell_1 < \ell_2 < \dots < \ell_s$
 216 then we also have $r_1 < r_2 < \dots < r_s$. This gives us a linear order on the intervals, and so we
 217 can maintain them in any balanced binary search tree. After adding the new interval $[\ell, r]$ to
 218 the collection, we can check if it is not contained in any of the already existing intervals, and
 219 if so find the already existing intervals that should be removed, with standard operations on the
 220 balanced binary search tree.

221 Now we explain how to implement **Step 1** and **Step 2** efficiently using T_y and the
 222 secondary structures associated with its nodes. Let i and (x_i, y_i) be the index and the
 223 threshold we are currently considering. We begin our discussion with **Step 2** and therefore
 224 assume that we already computed $\min_L(x_i, y_i)$ and $\max_R(x_i, y_i)$ for this threshold. Note
 225 that all thresholds (x_b, y_b) such that $x_i < x_b$ and $y_i < y_b$ have been already processed.
 226 Moreover, all thresholds (x_c, y_c) such that $x_i < x_c$ have been already processed and will not
 227 be considered in the future, so we don't need to be concerned with updating their answer.
 228 Hence, in **Step 2** we update all thresholds (x_b, y_b) such that $y_b < y_i$, regardless of the value
 229 of x_b . To this end, we consider every ancestor y_b of y_i such that $y_b < y_i$, plus the node y_i
 230 itself, and add the interval $[\min_L(x_i, y_i), \max_R(x_i, y_i)]$ to their corresponding T_{\min} and T_{\max} .
 231 To implement **Step 1**, we consider every ancestor y_c of y_i such that $y_c > y_i$, plus the node
 232 y_i itself, and we query their corresponding T_{\min} and T_{\max} . It can readily be verified that by the

choice of which ancestors are updated, this is enough to implicitly consider every $y_b > y_i$, as such y_b must have updated one of the ancestors y_c .

Time complexity. The time complexity of the algorithm is $\mathcal{O}(n \log^2 n)$. First, we need to sort the $4n$ thresholds in $\mathcal{O}(n \log n)$ time. Each of these thresholds is processed by considering $\mathcal{O}(\log n)$ nodes of T_y . At each of these nodes u we spend $\mathcal{O}(\log n)$ amortized time to update and query $T_{\min}(u)$ and $T_{\max}(u)$. Furthermore, the space complexity is $\mathcal{O}(n \log n)$, because each interval appears in $\mathcal{O}(\log n)$ secondary structures. Instead of using balanced binary search trees with $\mathcal{O}(\log n)$ query and update time for the secondary structures, we can plug in any predecessor structure that stores a collection of s integers from $\{1, 2, \dots, n\}$ in $\mathcal{O}(s)$ space with expected $\mathcal{O}(\log \log n)$ query and update time [33].

4 Order-Preserving Matching of Indeterminate Strings

Two determinate strings $p_w[1..n]$ and $t_w[1..n]$ are order-equivalent when, for every $i, j \in \{1, \dots, n\}$, $p_w[i] \leq p_w[j]$ iff $t_w[i] \leq t_w[j]$. This is denoted by $p_w \sim_{\leq} t_w$. Given two indeterminate strings p and t of equal-length n with at most 2 uncertain characters per position, we want to check if there exist $p_w \in \tilde{p}$ and $t_w \in \tilde{t}$ such that $p_w \sim_{\leq} t_w$. The goal of this section is to prove that this is NP-hard by reducing checking satisfiability of a 3-CNF formula.

We start with rephrasing the question in a graph-theoretical language. Let Σ_p and Σ_t be the sets of characters that occur in p and t , respectively. We consider a complete bipartite graph G with Σ_p corresponding to the nodes on the one side and Σ_t corresponding to the nodes on the other side. We claim that there exist $p_w \in \tilde{p}$ and $t_w \in \tilde{t}$ such that $p_w \sim_{\leq} t_w$ iff there exists a non-crossing matching M in G , where non-crossing means that we cannot have two edges $(x, y), (x', y')$ such that $x < x'$ but $y' < y$, such that the following holds for every position $i = 1, 2, \dots, n$:

$p[i] = x$ and $t[i] = y : (x, y) \in M$,
 $p[i] = x_1|x_2$ and $t[i] = y : (x_1, y) \in M$ or $(x_2, y) \in M$
 $p[i] = x$ and $t[i] = y_1|y_2 : (x, y_1) \in M$ or $(x, y_2) \in M$,
 $p[i] = x_1|x_2$ and $t[i] = y_1|y_2 : (x_1, y_1) \in M$ or $(x_1, y_2) \in M$ or $(x_2, y_1) \in M$ or $(x_2, y_2) \in M$.

The proof is straightforward.

We consider a 3-CNF formula ϕ on n variables $1, 2, \dots, n$ and m clauses. We reduce checking satisfiability of ϕ to finding a non-crossing matching M in a complete bipartite undirected graph G that respects a number of constraints of the form $M \cap X \times Y \neq \emptyset$, for $|X|, |Y| \leq 2$, or $X \times Y$ for short. As long as the size of G and the number of constraints is polynomial, this will establish NP-hardness of our problem, as we can create two strings p and t and encode each constraint by setting up some $p[i]$ and $t[i]$ appropriately.

We start with creating nodes $1, 2, \dots, n$ on the left side and $1, 2, 3, 4, \dots, 2n$ on the right side of G . We add a constraint $\{i\} \times \{2i-1, 2i\}$, for every $i = 1, 2, \dots, n$. We add a constraint $\{2n+1\} \times \{2n+1\}$. For every $k = 1, 2, \dots, m$, we consider the k -th clause $(\ell_{k,1} \vee \ell_{k,2} \vee \ell_{k,3})$, where $\ell_{k,1}, \ell_{k,2}, \ell_{k,3}$ are literals. Let $s = 2n + 2 + 5(k-1)$. We add the following constraints: $\{s\} \times \{s, s+1\}$, $\{s+2, s+3\} \times \{s+3\}$, $\{s, s+2\} \times \{s+1, s+3\}$. This is illustrated in Figure 5. Then we add a constraint for every literal:

1. If $\ell_{k,1} = x$ then we add $\{x, s\} \times \{2x, s\}$, and if $\ell_{k,1} = \neg x$ then we add $\{x, s\} \times \{2x-1, s\}$.

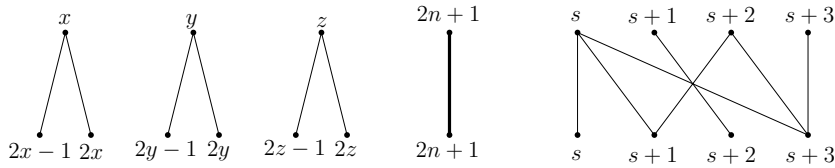
- 276 2. If $\ell_{k,2} = y$ then we add $\{y, s+1\} \times \{2y, s+2\}$, and if $\ell_{k,2} = \neg y$ then we add $\{x, s+1\} \times$
 277 $\{2y-1, s+2\}$.
 278 3. If $\ell_{k,3} = z$ then we add $\{z, s+3\} \times \{2z, s+3\}$, and if $\ell_{k,3} = \neg z$ then we add $\{z, s+3\} \times$
 279 $\{2z-1, s+3\}$.

280 Due to the constraint $\{2n+1\} \times \{2n+1\}$, a variable constraint $\{v, a\} \times \{2v-1, b\}$ translates
 281 into $(v, 2v-1) \in M$ or $(a, b) \in M$. Similarly, $\{v, a\} \times \{2v, b\}$ translates into $(v, 2v) \in M$ or
 282 $(a, b) \in M$.

283 We need to prove that ϕ is satisfiable iff there exists a non-crossing matching M in G
 284 that respects all the constraints.

285 First, assume that ϕ is satisfiable and fix a satisfying valuation of all the variables. We
 286 obtain M by first adding $(v, 2v-1)$ or $(v, 2v)$ to M depending on whether v is set to false
 287 or true, respectively. We also add $(2n+1, 2n+1)$ to M . Then, we proceed as follows for
 288 the k -th clause. For concreteness assume that the clause is $(x \vee y \vee z)$, the argument is
 289 symmetric for the other cases. If x is set to false then we add (s, s) to M . If y is set to false
 290 then we add $(s+1, s+2)$ to M . Finally, if z is set to false then we add $(s+3, s+3)$ to M .
 291 Because at least one of x, y, z is set to true, at least one of these three edges is not in M . If
 292 $(s+1, s+2) \notin M$ then we add $(s+2, s+1)$ to M . If $(s, s) \notin M$ then we add $(s, s+1)$ to M .
 293 Finally, if $(s+3, s+3) \notin M$ then we add $(s+2, s+3)$ to M . In all cases, the constraints
 294 corresponding to the k -clause are fulfilled. Due to how we compose the gadgets, M being
 295 a non-crossing matching in every gadget implies that M is a non-crossing matching in the
 296 whole G .

297 Second, assume that we have a non-crossing matching M in G . For every $v = 1, 2, \dots, n$,
 298 M contains exactly one of the edges $(v, 2v-1)$, $(v, 2v)$. We set v to false if $(v, 2v-1) \in M$
 299 and to true if $(v, 2v) \in M$. We must have $(2n+1, 2n+1) \in M$. We need to verify that every
 300 clause is satisfied by the obtain valuation of the variables. Again, for concreteness assume
 301 that the clause is $(x \vee y \vee z)$. We cannot have all edges (s, s) , $(s+1, s+2)$, $(s+3, s+3)$ in
 302 M , as in such case the constraint $\{s, s+2\} \times \{s+1, s+3\}$ cannot be fulfilled. If $(s, s) \notin M$
 303 then due to the constraint $\{x, s\} \times \{2x, s\}$ we must have $(x, 2x) \in M$, so x is set to true.
 304 If $(s+1, s+2) \notin M$ then due to the constraint $\{y, s+1\} \times \{2y, s+2\}$ we must have
 305 $(y, 2y) \in M$, so y is set to true. Finally, if $(s+3, s+3) \notin M$ then due to the constraint
 306 $\{z, s+3\} \times \{2z, s+3\}$ we must have $(z, 2z) \in M$, so z is set to true. So, one of the variable
 307 x, y, z is set to true, making the clause satisfied.



■ **Figure 5** Gadget created for the k -th clause concerning variables x, y, z .

308 5 Parameterized Matching of Indeterminate Strings

309 Two determinate strings $p_w[1..n]$ and $t_w[1..n]$ are parameterized-equivalent when, for every
 310 $i, j \in \{1, \dots, n\}$, $p_w[i] = p_w[j]$ iff $t_w[i] = t_w[j]$. This is denoted by $p_w \sim_w t_w$. Given two
 311 indeterminate strings p and t of equal-length n with at most 2 uncertain characters per
 312 position, we want to check if there exist $p_w \in \tilde{p}$ and $t_w \in \tilde{t}$ such that $p_w \sim_w t_w$. The goal of

313 this section is to prove that this is NP-hard by reducing checking if a given undirected graph
 314 has a vertex cover consisting of at most k vertices.

315 As in the previous section, we start with rephrasing the question in a graph-theoretical
 316 language. Let Σ_p and Σ_t be the sets of characters that occur in p and t , respectively. We
 317 consider a complete bipartite graph G with Σ_p corresponding to the nodes on the one side
 318 and Σ_t corresponding to the nodes on the other side. We claim that there exist $p_w \in \tilde{p}$ and
 319 $t_w \in \tilde{t}$ such that $p_w \sim_t t_w$ iff there exists a matching M in G , such that the following holds
 320 for every position $i = 1, 2, \dots, n$:

321 $p[i] = x$ and $t[i] = y : (x, y) \in M$,
 322 $p[i] = x_1|x_2$ and $t[i] = y : (x_1, y) \in M$ or $(x_2, y) \in M$
 323 $p[i] = x$ and $t[i] = y_1|y_2 : (x, y_1) \in M$ or $(x, y_2) \in M$,
 324 $p[i] = x_1|x_2$ and $t[i] = y_1|y_2 : (x_1, y_1) \in M$ or $(x_1, y_2) \in M$ or $(x_2, y_1) \in M$ or
 325 $(x_2, y_2) \in M$.

326 The proof is straightforward.

327 We consider an undirected graph H on n vertices $V = \{1, 2, \dots, n\}$ and m edges E
 328 together with a parameter $k \leq n$. We reduce checking if there is a subset S of k vertices of
 329 H such that for every edge $(u, v) \in E$ we have $u \in S$ or $v \in S$ to finding a matching M in a
 330 complete bipartite graph G that respects a number of constraints of the form $M \cap X \times Y \neq \emptyset$,
 331 for $|X|, |Y| \leq 2$, or $X \times Y$ for short. As long as the size of G and the number of constraints
 332 is polynomial, this will establish NP-hardness of our problem, as we can create two strings p
 333 and t and encode each constraint by setting up some $p[i]$ and $t[i]$ appropriately.

334 We start with creating nodes $1, 2, \dots, n$ on the left side and $1, 2, \dots, n$ on the right side
 335 of G . We add a constraint $\{u, v\} \times \{u, v\}$ for every $(u, v) \in E$. For every $i = 1, 2, \dots, n$,
 336 $(i, j) \in M$ for some $j \in \{1, 2, \dots, n\}$ corresponds to including i in the sought vertex cover.
 337 The remaining part of H is constructed as to guarantee that there are at least k nodes
 338 $i \in \{1, 2, \dots, n\}$ such that $(i, j) \in M$ for some $j \neq \{1, 2, \dots, n\}$. To this end, we design a
 339 gadget G_{2s} with the following property:

- 340 1. there are distinguished $2s$ nodes v_1, v_2, \dots, v_{2s} on the left side, each v_i is incident to a
 341 unique edge e_i ,
- 342 2. there are also some additional internal nodes on the left and on the right and some
 343 constraints that concern both the internal and the distinguished nodes,
- 344 3. if none of the edges e_i belongs to M then it is not possible to satisfy the constraints of
 345 G_{2s} ,
- 346 4. for any nonempty subset S of distinguished nodes, it is possible to select some of the
 347 edges with both endpoints being internal nodes in such a way that, together with the
 348 edges e_i for $i \in S$, they satisfy all constraints of G_{2s} .

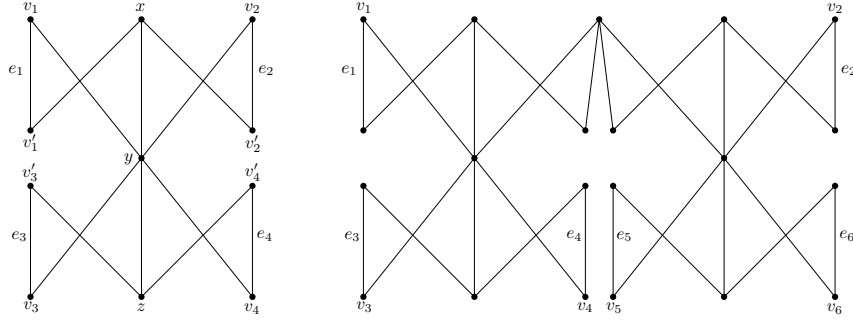
349 We will first show that G_4 exists, and then explain how to obtain $G_{2(s+1)}$ from G_{2s} .

350 ► **Lemma 7.** G_4 with the sought properties exists.

351 **Proof.** G_4 consists of nodes v_1, v_2, v_3, v_4 and internal nodes v'_1, v'_2, v'_3, v'_4 and x, y, z . We
 352 set $e_i = (v_i, v'_i)$ for $i = 1, 2, 3, 4$ and create the following constraints: $\{v_1, x\} \times \{v'_1, y\}$,
 353 $\{v_2, x\} \times \{v'_2, y\}$, $\{v_3, z\} \times \{v'_3, y\}$, $\{v_4, z\} \times \{v'_4, y\}$ and $\{y\} \times \{x, z\}$. See Figure 6.

354 Assume that none of the edges e_i belongs to M . By symmetry, we can assume that
 355 $(x, y) \in M$. But then we must have $(v'_3, z), (v'_4, z) \in M$, which is impossible.

356 Let S be a nonempty set of distinguished nodes. By symmetry, we can assume that
 357 $v_1 \in S$. Then, we include $(y, z) \in M$ and if $e_2 \notin S$ we also include $(v'_2, x) \in M$. ◀



■ **Figure 6** Gadgets G_4 (left) and G_8 (right).

358 ► **Lemma 8.** $G_{2(s+1)}$ with the sought properties can be obtained in polynomial time from
 359 G_{2s} with the sought properties.

360 **Proof.** We take a copy of G_{2s} , let its distinguished nodes and their corresponding edges
 361 be v_1, v_2, \dots, v_{2s} and e_1, e_2, \dots, e_{2s} . We also take a copy of G_4 , let its distinguished nodes
 362 and their corresponding edges be u_1, u_2, u_3, u_4 and f_1, f_2, f_3, f_4 . To obtain $G_{2(s+1)}$ we
 363 identify v_{2s} with u_1 and add a constraint that enforces including e_{2s} or f_1 in M . The
 364 distinguished nodes and their corresponding edges of $G_{2(s+1)}$ are $v_1, v_2, \dots, v_{2s-1}, u_2, u_3, u_4$
 365 and $e_1, e_2, \dots, e_{2s-1}, f_2, f_3, f_4$.

366 Assume that none of the edges $e_1, e_2, \dots, e_{2s-1}, f_2, f_3, f_4$ belongs to M . Either $e_{2s} \notin M$
 367 and we obtain that none of the edges e_1, e_2, \dots, e_{2s} belongs to M , or $f_1 \notin M$ and none of the
 368 edges f_1, f_2, f_3, f_4 belongs to M . In either case we obtain a contradiction by the construction
 369 of G_{2s} or G_4 .

370 Let S be a nonempty set of distinguished nodes and assume that $v_1 \in S$ (other cases are
 371 essentially the same). We set $S' = S \cap \{v_1, v_2, \dots, v_{2s-1}\}$ and $S'' = (S \cap \{u_2, u_3, u_4\}) \cup \{u_1\}$.
 372 Then $S', S'' \neq \emptyset$, and by assumption we can select some of the edges with both endpoints
 373 being internal nodes of G_{2s} or G_4 in such a way that, together with the edges e_i for $i \in S'$
 374 and f_j for $j \in S''$, they satisfy all constraints of G_{2s} and G_4 . Additionally, the constraint
 375 that enforces including e_{2s} or f_1 is satisfied by taking f_1 . So, by selecting the edges with
 376 both endpoints being internal nodes of G_{2s} or G_4 together with f_1 we obtain a set of edges
 377 with both endpoints being internal nodes of $G_{2(s+1)}$ that, together with the edges associated
 378 with the nodes in S , satisfy all constraints of $G_{2(s+1)}$ as required. ◀

379 With the gadget G_{2s} in hand, we are ready to complete the reduction. By duplicating
 380 the graph H we can assume that $n = 2s$. We add $n - k$ copies of the gadget G_{2s} to G . Let
 381 v_1, v_2, \dots, v_{2s} be the distinguished nodes of one such copy. We identify v_i with the node i
 382 on the left side of G . This guarantees that for each gadget we must have a unique node i
 383 such that $(i, 1), (i, 2), \dots, (i, n) \notin M$. We claim that the resulting graph G has a matching
 384 that satisfies all the constraints if and only if H admits a vertex cover of cardinality at most
 385 k . In one direction, consider the set C consisting of all nodes $i \in \{1, 2, \dots, n\}$ such that
 386 $(i, 1), (i, 2), \dots, (i, n) \notin M$. By the properties of G_{2s} , $|C| \leq k$. We need to argue that C is a
 387 vertex cover. Consider any $(u, v) \in E$. Due to the constraint $\{u, v\} \times \{u, v\}$, one of the edges
 388 $(u, u), (u, v), (v, u), (v, v)$ must belong to M . But then either u or v cannot be matched to
 389 any node not belonging to $\{1, 2, \dots, n\}$, so $u \in C$ or $v \in C$ as required. In other direction,
 390 let C be a vertex cover of H of cardinality at most k . For every $i \in C$, we include the edge
 391 (i, i) in M . This clearly satisfies every constraint $\{u, v\} \times \{u, v\}$ by C being a vertex cover.

Then, for every copy of G_{2s} we choose a unique node $i \notin C$ (that is not matched to any other node yet) and use the properties of G_{2s} to add its internal edges to M in such a way that, together with the edge associated to i , they satisfy all the constraints.

References

- 1 A. Alatabbi, A. S. S. Islam, M. S. Rahman, J. Simpson, and W. F. Smyth. Enhanced covers of regular & indeterminate strings using prefix tables, 2015. [arXiv:1506.06793](#).
- 2 A. Amir, Y. Aumann, G. M. Landaud, M. Lewenstein, and N. Lewenstein. Pattern matching with swaps. *Journal of Algorithms*, 37(2):247–266, 2000.
- 3 A. Amir, M. Farach, and S. Muthukrishnan. Alphabet dependence in parameterized matching. *Inf. Proc. Lett*, 49(3):111–115, 1994.
- 4 Lewenstein M. Apostolico A., Erdős Péter L. Parameterized matching with mismatches. *J. Discrete Algorithms*, 5(1):135–140, 2007.
- 5 B. Baker. A theory of parameterized pattern matching: algorithms and applications. *STOC*, pages 71–80, 1993.
- 6 B. S. Baker. Parameterized pattern matching: Algorithms and applications. *Journal of Computer and System Sciences*, 52(1):28–42, 1996.
- 7 M. Bataa, S. G. Park, A. Amir, G. M. Landau, and K. Park. Finding periods in cartesian tree matching. *Colbourn C., Grossi R., Pisanti N. (eds) Combinatorial Algorithms. IWOC2019*, 11638:70–84, 2019.
- 8 G. Bernardini, P. Gawrychowski, N. Pisanti, S. P. Pissis, and G. Rosone. Even faster elastic-degenerate string matching via fast matrix multiplication, 2019. [arXiv:1905.02298](#).
- 9 P. Bille, I. Gørtz, H. Vildhøj, and D. Wind. String matching with variable length gaps. *Theoretical Computer Science*, 443:385–394, 10 2010.
- 10 P. Burcsi, F. Cicalese, G. Fici, and Z. Liptak. Algorithms for jumbled pattern matching in strings. *International Journal of Foundations of Computer Science*, 23(02):357–374, 2012.
- 11 S. Choa, J. C. Na, K. Park, and J. S. Sima. A fast algorithm for order-preserving pattern matching. *Information Processing Letters*, 115:397–402, 2015.
- 12 M. Christodoulakis, P.J. Ryan, W.F. Smyth, and S. Wang. Indeterminate strings, prefix arrays & undirected graphs. *Theoretical Computer Science*, 600:34–48, 2015.
- 13 D. Costa, L. M. S. Russo, R. Henriques, H. Bannai, and A. P. Francisco. Order-preserving pattern matching indeterminate strings. 2019. [arXiv:1905.02589](#).
- 14 M. Crochemore, C. S. Iliopoulos, T. Kociumaka, J. Radoszewski, W. Rytter, and T. Waleń. Covering problems for partial words and for indeterminate strings, 2014. [arXiv:1412.3696](#).
- 15 J. W. Daykin, R. Groult, Y. Guesnet, T. Lecroq, A. Lefebvre, M. Léonard, L. Mouchard, É. Prieur-Gaston, and B. Watson. Efficient pattern matching in degenerate strings with the burrows-wheeler transform, 2017. [arXiv:1708.01130](#).
- 16 P. Gawrychowski and P. Uznański. Order-preserving pattern matching with k mismatches. *Theoretical Computer Science*, 638:136–144, 2016. [doi:10.1016/j.tcs.2015.08.022](#).
- 17 G. Gourdel, T. Kociumaka, J. Radoszewski, W. Rytter, A. Shur, and T. Waleń. String periods in the order-preserving model. 96(38):1–16, 2018.
- 18 G. Gu, S. Song, S. Faro, T. Lecroq, and K. Park. Fast multiple pattern cartesian tree matching. *WALCOM2020*, 2019.
- 19 C. Hazay, M. Lewenstein, and D. Sokol. Approximate parameterized matching. *ACM Transactions on Algorithms (TALG)*, 3(3):29–44, 2007.
- 20 J. Helling, P. Ryan, W.F. Smyth, and M. Soltys. Constructing an indeterminate string from its associated graph. *Theoretical Computer Science*, 710, 03 2017. [doi:10.1016/j.tcs.2017.02.016](#).
- 21 Rui Henriques, Alexandre P. Francisco, Luís M. S. Russo, and Hideo Bannai. Order-preserving pattern matching indeterminate strings. In *CPM*, volume 105 of *LIPIcs*, pages 2:1–2:15. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2018.

- 442 22 J. Holub, W. F. Smyth, and S. Wang. Fast pattern-matching on indeterminate strings. *Discrete*
443 *Algorithms*, 6(1):37–50, 2008.
- 444 23 J. Holub and W.F. Smyth. Algorithms on indeterminate strings. in: *M. Miller, K. Park*
445 *(Eds.), Proc. 14th Australasian Workshop on Combinatorial Algorithms*, pages 36 – 45, 2003.
- 446 24 C. Iliopoulos, R. Kundu, and S. Pissis. Efficient pattern matching in elastic-degenerate strings,
447 2016. [arXiv:1610.08111](#).
- 448 25 J. Kima, P. Eades, R.Fleischer, S. H. Hong, C. S. Iliopoulou, K. Park, S. J. Puglisig, and
449 T. Tokuyamah. Order-preserving matching. *Theoretical Computer Science*, 525:68–79, 2014.
- 450 26 M. Kubica, T. Kulczyński, J. Radoszewski, W. Rytter, and T. Waleńca. A linear time
451 algorithm for consecutive permutation pattern matching. *Information Processing Letters*,
452 113(12):430–433, 2013. doi:10.1016/j.ipl.2013.03.015.
- 453 27 R. McIntyre and M. Soltys. An exact upper bound on the size of minimal clique covers, 2017.
454 [arXiv:1705.06326](#).
- 455 28 S. Park, A. Amir, G. M. Landau, and K. Park. Cartesian tree matching and indexing. *CPM*,
456 128, 2019.
- 457 29 M. S. Rahman and C. S. Iliopoulos. Pattern matching algorithms with don’t cares. 2007.
- 458 30 S. Song, C. Ryu, S. Faro, T. Lecroq, and K. Park. Fast cartesian tree matching. *SPIRE*, 2019.
- 459 31 J. Vuillemin. A unifying look at data structures. *Communications of the ACM*, 23(4):229–239,
460 1980.
- 461 32 R. A. Wagner and M. J. Fischer. The string-to-string correction problem. *J. ACM*,
462 21(1):168–173, 1974.
- 463 33 Dan E. Willard. Log-logarithmic worst-case range queries are possible in space $\Theta(N)$. *Inf.*
464 *Process. Lett.*, 17(2):81–84, 1983.
- 465 34 B. Zeidman. Software v. software. *IEEE Spectrum*, 47:32–53, 2010.